

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
30 January 2003 (30.01.2003)

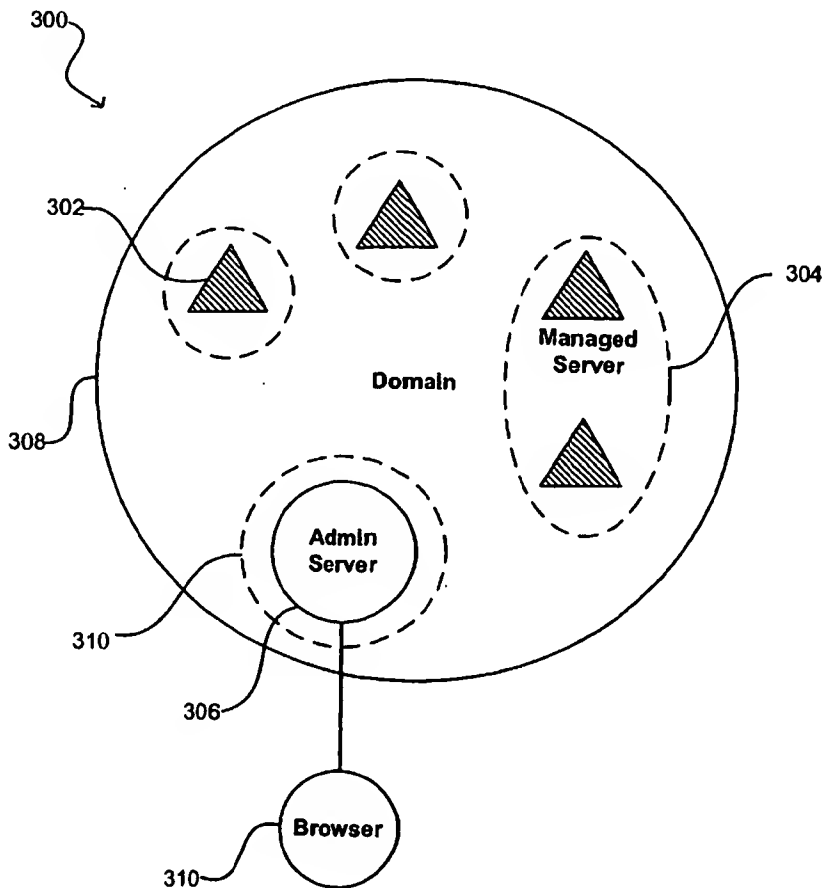
PCT

(10) International Publication Number
WO 03/009092 A2

- (51) International Patent Classification⁷: **G06F** (72) Inventors: **JACOBS, Dean, Bernard**; 1747 Madera Street, Berkeley, CA 94707 (US). **KRAMER, Reto**; 411 Green Street, #2A, San Francisco, CA 94133 (US). **SRINIVASAN, Ananthan, Bala**; 1610 Sanchez Street, San Francisco, CA 94131 (US).
- (21) International Application Number: PCT/US02/22366
- (22) International Filing Date: 15 July 2002 (15.07.2002)
- (25) Filing Language: English (74) Agents: **MEYER, Sheldon, R. et al.**; Fliesler Dubb Meyer & Lovejoy LLP, Four Embarcadero Center, Fourth Floor, San Francisco, CA 94111-4156 (US).
- (26) Publication Language: English
- (30) Priority Data:
- | | | |
|------------|------------------------------|----|
| 60/305,986 | 16 July 2001 (16.07.2001) | US |
| 607305,978 | 16 July 2001 (16.07.2001) | US |
| 09/975,590 | 11 October 2001 (11.10.2001) | US |
| 09/975,587 | 11 October 2001 (11.10.2001) | US |
- (71) Applicant: **BEA SYSTEMS, INC.** [US/US]; 2315 North First Street, San Jose, CA 95131 (US).
- (81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZM, ZW.

[Continued on next page]

(54) Title: DATA REPLICATION PROTOCOL



(57) Abstract: Data can be replicated over a network using a one or two phase method. For the one phase method, a master server containing an original copy of the data sends a version number for the current state of the data to each slave on the network so that each slave can request a delta from the master. The delta that is requested contains the data necessary to update the slave to the appropriate version of the data. For the two phase method, the master server sends a packet of information to each slave. The packet of information can be committed by the slaves if each slave is able to process the commit.

WO 03/009092 A2



(84) Designated States (regional): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

— without international search report and to be republished upon receipt of that report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

DATA REPLICATION PROTOCOL

CLAIM OF PRIORITY

5 **[0001]** This application claims priority to U.S. Provisional patent application No. 60/305,986, filed July 16, 2001, entitled DATA REPLICATION PROTOCOL; U.S. Provisional patent application No. 60/305,978, filed July 16, 2001 entitled LAYERED ARCHITECTURE FOR DATA REPLICATION; U.S. Patent application No. 09/975,590, filed 10 October 11, 2001, entitled DATA REPLICATION PROTOCOL; U.S. Patent application No. 09/975,587, filed October 11, 2001 entitled LAYERED ARCHITECTURE FOR DATA REPLICATION incorporated herein by reference.

COPYRIGHT NOTICE

15 **[0004]** A portion of the disclosure of this patent document contains material which is subject to copyright protection. The copyright owner has no objection to the facsimile reproduction by anyone of the patent document or the patent disclosure, as it appears in the Patent and 20 Trademark Office patent file or records, but otherwise reserves all copyright rights whatsoever.

TECHNICAL FIELD

25 **[0005]** The invention relates generally to a system for transferring data. The invention relates more specifically to a system and method for replicating data over a network.

BACKGROUND

30 **[0006]** There are several types of distributed processing systems. Generally, a distributed processing system includes a plurality of processing devices, such as two computers coupled through a communication medium. One type of distributed processing system is a

client/server network. A client/server network includes at least two processing devices, typically a central server and a client. Additional clients may be coupled to the central server, there may be multiple servers, or the network may include only servers coupled through the communication medium.

[0007] In such a network environment, it is often desirable to send applications or information from the central server to a number of workstations and/or other servers. Often, this may involve separate installations on each workstation, or may involve separately pushing a new library of information from the central server to each individual workstation and/or server. These approaches can be time consuming and are an inefficient use of resources. The separate installation of applications on each workstation or server also introduces additional potential sources of error.

[0008] Ideally, the sending of information should be both reliable in the face of failures and scalable, so that the process makes efficient use of the network. Conventional solutions generally fail to achieve one or both of these goals. One simple approach is to have a master server individually contact each slave and transfer the data over a point-to-point link, such as a TCP/IP connection. This approach leads to inconsistent copies of the data if one or more slaves are temporarily unreachable, or if the slaves encounter an error in processing the update. At the other extreme are complex distributed agreement protocols, which require considerable cross-talk among the slaves to ensure that all copies of the data are consistent.

BRIEF SUMMARY

[0009] The present invention includes a method for replicating data from a master server to at least one slave or managed server, such as may be accomplished on a network. In the method, it may be determined whether the replication should be accomplished in a one or two phase

method. If the replication is to be accomplished in a one phase method, a version number may be sent that corresponds to the current state of the data on the master server. This version number may be sent to every slave server on the network, or only a subset of slave servers. The slave servers
5 receiving the version number may then request that a delta be sent from the master. The delta may contain data necessary to update the data on that slave to correspond to the current version number.

[0010] If the replication is to be accomplished in a two phase method, a packet of information may be sent from the master to each
10 slave, or a subset of slaves. Those slaves may then respond to the master server whether they can commit the packet of information. If at least some of the slaves can commit the data, the master may signal to those slave that they should process the commit. After processing the commit, those slaves may update to the current version number. If any of the slaves are
15 unable to process the commit, the commit may be aborted.

BRIEF DESCRIPTION OF THE DRAWINGS

[0011] Figure 1 is a diagram of a domain structure in accordance with one embodiment of the present invention.

20 [0012] Figure 2 is a diagram of layered architecture in accordance with one embodiment of the present invention.

[0013] Figure 3 is a diagram of a clustered domain structure in accordance with one embodiment of the present invention.

25 [0014] Figure 4 is a diagram of one phase process for a layered architecture in accordance with one embodiment of the present invention.

[0015] Figure 5 is a diagram of two phase process for a layered architecture in accordance with one embodiment of the present invention.

[0016] Figure 6 is a flowchart for a one phase process in accordance with one embodiment of the present invention.

30 [0017] Figure 7 is a flowchart for a two phase process in accordance with one embodiment of the present invention.

DETAILED DESCRIPTION

[0018] The present invention provides for the replication of data or other information, such as from a master server, or "administration" server ("Admin server"), to a collection of slave servers, or "managed" servers.

5 This replication can occur over any appropriate network, such as a conventional local area network or ethernet. In one embodiment, a master server owns the original record of all data on the network, to which any updates are to be applied. A copy of the data, together with updates as they occur, can be transmitted to each slave server. One example
10 application involves the distribution of configuration information from an Admin server to a collection of managed servers.

[0019] In one system in accordance with the present invention, it may be necessary for a service, such as a Data Replication Service (DRS), to distribute configuration and deployment information from an Admin
15 Server to managed servers in the appropriate domain. Large data items can be distributed over point-to-point connections, such as Transmission Control Protocol ("TCP"), since a multicast protocol like User Datagram Protocol ("UDP") does not have flow control, and can overwhelm the system. Remote Method Invocation (RMI), Hypertext Transfer Protocol
20 (HTTP), or a similar protocol may be used for point-to-point connections.

[0020] Managed servers can also persistently cache data on local disks. Without such caching, an unacceptable amount of time may be required to transfer the necessary data. The ability of the managed servers to cache is important, as it increases the speed of startup by reducing the
25 amount of startup data to be transferred. Caching can also allow startup and/or restart if the Admin Server is unreachable. Restart may be a more attractive option, and it may be the case that the Admin server directs a server to start. Caching, however, can provide the ability to start the domain without the Admin Server being available.

30 **[0021]** As shown in the domain structure 100 of **Figure 1**, an Admin Server 102 and at least one managed server 104 can comprise a domain

106. This domain 106 can be the administration unit for startup and shutdown. In one embodiment, a browser 108, or other user application or device, tells the Admin Server 102 to start. The Admin Server 102 then tells all managed servers 104 in the domain 106 to start, and passes the
5 appropriate configuration information. If a server goes down after the managed servers 104 have started, it may be desirable for that server to restart automatically, whether or not the Admin Server 102 is available. Cached data can be useful for this purpose.

[0022] Updates to data on the Admin Server can be packaged as
10 incremental deltas between versions. The deltas can contain configuration and/or other information to be changed. It may be preferable to update the configuration while the domain is running, as it may be undesirable to take the system offline. In one embodiment, the configuration changes happen dynamically, as they are pushed out by the Admin Server. Only the
15 changes to the configuration are sent in the deltas, as it may be unnecessary, and unduly cumbersome, to send the full configuration each time.

[0023] A protocol in accordance with the present invention integrates two methods for the distribution of updates, although other appropriate
20 methods may be used accordingly. These distribution methods may be referred to as a one-phase method and a two-phase method, and can provide a tradeoff between consistency and scalability. In a one-phase method, which may favor scalability, each slave can obtain and process updates at its own pace. Slaves can get updates from the master at
25 different times, but can commit to the data as soon as it is received. A slave can encounter an error in processing an update, but in the one-phase method this does not prevent other slaves from processing the update.

[0024] In a two-phase method in accordance with the present invention, which may favor consistency, the distribution can be "atomic", in
30 that either all or none of the slaves successfully process the data. There can be separate phases, such as prepare and commit phases, which can

allow for a possibility of abort. In the prepare phase, the master can determine whether each slave can take the update. If all slaves indicate that they can accept the update, the new data can be sent to the slaves to be committed in the commit phase. If at least one of the slave servers cannot take the update, the update can be aborted and there may not be a commit. In this case, the managed servers can be informed that they should roll back the prepare and nothing is changed. Such a protocol in accordance with the present invention is reliable, as a slave that is unreachable when an update is committed, in either method, eventually gets the update.

[0025] A system in accordance with the present invention can also ensure that a temporarily unavailable server eventually receives all updates. For example, a server may be temporarily isolated from the network, then come back into the network without restarting. Since the server is not restarting, it normally will not check for updates. The server coming back into the network can be accounted for by having the server check periodically for new updates, or by having a master server check periodically to see whether the servers have received the updates.

[0026] In one embodiment, a master server regularly sends multicast "heartbeats" to the slave servers. Since a multicast approach can be unreliable, it is possible for a slave to miss arbitrary sequences of heartbeats. For instance, a slave server might be temporarily disconnected from the network due to a network partitioning, or the slave server itself might be temporarily unavailable to the network due, causing a heartbeat to be missed. Heartbeats can therefore contain a window of information about recent updates. Such information about previous updates may be used to reduce the amount of network traffic, as explained below.

[0027] There can be at least two layers within each master and each slave: a user layer and a system layer (or DRS layer). The user layer can correspond to the user of the data replication system. A DRS layer can

correspond to the implementation of the data replication system itself. The interaction of these participants and layers is shown in **Figure 2**.

[0028] As shown in the startup diagram 200 of **Figure 2**, the master user 202 and slave user 204 layers in this embodiment make downcalls into the master DRS 206 and slave DRS 208 layers, respectively. Such
5 downcalls can, for example, take the form of:

registerMaster(DID, verNum, listener)

registerSlave(DID, verNum, listener)

where *DID* is an identifier taken from knowledge of well-known DIDs and
10 refers to the object of interest, *verNum* is taken from the local persistent store as the user's current version number, and *listener* is an object that will handle upcalls from the DRS layer. The upcall can call a method on the listener object. The master can then begin to send heartbeats, or periodic deltas, with the current version number. A container layer 210 is shown,
15 which can include containers adapted to take information from the slave user 204. Examples of possible containers include enterprise Java beans, web interfaces, and J2EE (Java 2 Platform, Enterprise Edition) applications. Other applications and/or components can plug into the container layer 210, such as an administration client 212. Examples of
20 update messaging between the User and DRS layers are shown for the one phase method in **Figure 4**, as well as for the two phase method in **Figure 5**.

[0029] **Figure 4** shows one basic process 400 that may be used for a one-phase distribution approach in a layered architecture in accordance
25 with the present invention. In this process, the master user layer 402 makes a downcall 404 into the master DRS layer 406 to start a one phase distribution. This call can be to all slaves in the system, or only to a subset of slave servers. If the call is to a subset, the master user layer 402 can determine the scope of the update, or which slaves should receive the
30 update.

[0030] The master DRS layer begins multicasting heartbeats **408**, containing the current version number of the data on the master, to the slave DRS layer **410**. The slave DRS layer **410** requests the current version number **412** for the slave from the slave user layer **414**. The slave user layer **414** then responds **416** to the slave DRS layer **416** with the slave version number. If the slave is in sync, or already is on the current version number, then no further requests may be made until the next update. If the slave is out-of-sync and the slave is in the scope of the update, the slave DRS layer **410** can request a delta **420** from the master DRS layer **406** in order to update the slave to the current version number of the data on the master. The master DRS layer **406** requests **422** that the master user layer **402** create a delta to update the slave. The master user layer **402** then sends the delta **424** to the master DRS layer **406**, which forwards the delta **426** and the current version number of the master to the slave DRS layer **410**, which sends the delta **426** to the slave user to be committed. The current version number is sent with the delta in case the master has updated since the heartbeat **408** was received by the slave.

[0031] The master DRS layer **406** can continue to periodically send a multicast heartbeat containing the version number **408** to the slave server(s). This allows any slave that was unavailable, or unable to receive and process a delta, to determine that it is not on the current version of the data and request a delta **420** at a later time, such as when the slave comes back into the system.

[0032] Figure 5 shows one basic process **500** that may be used for a two phase distribution approach in a layered architecture in accordance with the present invention. In this process, the master user layer **504** makes a downcall **504** into the master DRS layer **506** to start a two phase distribution. The master user layer **502** may again need to determine the scope of the update, and may set a "timeout" value for the update process.

[0033] The master DRS layer **506** sends the new delta **508** to the

slave DRS layer 510. The slave DRS layer 510 sends a prepare request 512 to the slave user layer 514 for the new delta. The slave user layer 514 then responds 516 to the slave DRS layer 510 whether or not the slave can process the new delta. The slave DRS layer forwards the response 518 to the master DRS layer 506. If the slave cannot process the request because it is out-of-sync, the master DRS layer 506 makes an upcall 520 to the master user layer 502 to create a delta that will bring the slave in sync to commit the delta. The master user layer 502 sends the syncing delta 522 to the master DRS layer, which forwards the syncing delta 524 to the slave DRS layer 510. If the slave is able to process the syncing delta, the slave DRS layer 510 will send a sync response 526 to the master DRS layer 506 that the slave can now process the new delta. If the slave is not able to process the syncing delta, the slave DRS layer 510 will send the appropriate sync response 526 to the master DRS layer 506. The master DRS layer 506 then heartbeats a commit or abort message 528 to the slave DRS layer 510, depending on whether or not the slave responded that it was able to process the new delta. If all slave were able to prepare the delta, for example, the master can heartbeat a commit signal. Otherwise, the master can heartbeat an abort signal. The heartbeats also contains the scope of the update, such that a slave knows whether or not it should process the information contained in the heartbeat.

[0034] The slave DRS layer forwards this command 530 to the slave user layer 514, which then commits or aborts the update for the new delta. If the prepare phase was not completed within a timeout value set by the master user layer 502, the master DRS layer 506 can automatically heartbeat an abort 528 to all the slaves. This may occur, for example, when the master DRS layer 506 is unable to contact at least one of the slaves to determine whether that slave is able to process the commit. The timeout value can be set such that the master DRS layer 506 will try to contact the slave for a specified period of time before aborting the update.

[0035] For an update in a one-phase method, these heartbeats can cause each slave to request a delta starting from the slave's current version of the data. Such a process is shown in the flowchart of **Figure 6**. In this basic process **600**, which may or may not utilize a layered architecture in accordance with the present invention, a version number for the current data on the master server is sent from a master server to a slave server **602**. The slave server determines whether it has been updated to the current version number **604**. If the slave is not on the current version, it will request that a delta be sent from the master server to update the slave server **606**. When the delta is sent to the slave server, the slave server will process the delta in order to update the slave data to the current version **608**. The slave server will then update its version number to the current version number **610**.

[0036] For an update in a two-phase method, the master can begin with a prepare phase in which it pro-actively sends each slave a delta from the immediately-previous version. Such a process is shown in the flowchart of **Figure 7**. In this basic process **700**, which may or may not utilize a layered architecture in accordance with the present invention, a packet of information is sent from the master to a slave server or slave servers **702**. Each slave server receiving the packet determines whether it can process that packet and update to the current version **704**. Each slave server receiving the packet responds to the master server, indicating whether the slave server can process the packet **706**. If all slaves (to which the delta is sent) acknowledge successful processing of the delta within some timeout period, the master may decide to commit the update. Otherwise, the master server may decide to abort the update. Once this decision is made, the master server sends a message to the slave server(s) indicating whether the update should be committed or aborted **708**. If the decision is to commit, each server processes the commit **710**. Heartbeats may further

be used to signal whether a commit or abort occurred, in case the command was missed by one of the slaves.

[0037] A slave can be configured to immediately start and/or restart using cached data, without first getting the current version number from the master. As mentioned above, one protocol in accordance with the present invention allows slaves to persistently cache data on local disks. This caching decreases the time needed for system startup, and improves scalability by reducing the amount of data needing to be transferred. The protocol can improve reliability by allowing slaves to startup and/or restart if the master is unreachable, and may further allow updates to be packaged as incremental deltas between versions. If no cache data exists, the slave can wait for the master or can pull the data itself. If the slave has the cache, it may still not want to start out of sync. Startup time may be decreased if the slave knows to wait.

[0038] The protocol can be bilateral, in that a master or slave can take the initiative to transfer data, depending upon the circumstances. For example, a slave can pull a delta from the master during domain startup. When the slave determines it is on a different version than the delta is intended to update, the slave can request a delta from its current version to the current system version. A slave can also pull a delta during one-phase distribution. Here, the system can read the heartbeat, determine that it has missed the update, and request the appropriate delta.

[0039] A slave can also pull a delta when needed to recover from exceptional circumstances. Exceptional circumstances can exist, for example, when components of the system are out of sync. When a slave pulls a delta, the delta can be between arbitrary versions of the data. In other words, the delta can be between the current version of the slave and the current version of the system (or domain), no matter how many iterations apart those versions might be. In this embodiment, the availability of a heartbeat and the ability to receive deltas can provide synchronization of the system.

[0040] In addition to the ability of a slave to pull a delta, a master can have the ability to push a delta to a slave during two-phase distribution. In one embodiment, these deltas are always between successive versions of the data. This two-phase distribution method can minimize the likelihood of inconsistencies between participants. Slave users can process a prepare as far as possible without exposing the update to clients or making the update impossible to roll back. This can include such tasks as checking the servers for conflicts. If any of the slaves signals an error, such as by sending a "disk full" or "inconsistent configuration" message, the update can be uniformly rolled back.

[0041] It is still possible, however, that inconsistencies may arise. For instance, there may be errors in processing a commit, for reasons such as an inability to open a socket. Servers can also commit and expose the update at different times. Because the data cannot reach every managed server at exactly the same time, there can be some rippling effect. The use of multicasting can provide for a small time window, in an attempt to minimize the rippling effect. In one embodiment, a prepared slave will abort if it misses a commit, whether it missed the signal, the master crashed, etc.

[0042] A best-effort approach to multicasting can cause a slave server to miss a commit signal. If a master crashes part way through the commit phase, there may be no logging or means for recovery. There may be no way for the master to tell the remaining slaves that they need to commit. Upon abort some slaves may end up committing the data if the version is not properly rolled back. In one embodiment, the remaining slaves could get the update using one-phase distribution. This might happen, for example, when a managed server pulls a delta in response to a heartbeat received from an Admin server. This approach may maintain system scalability, which might be lost if the system tied down distribution in order to avoid any commit or version errors.

[0043] Each data item managed by the system can be structured to have a unique, long-lived domain identifier (DID) that is well-known across

the domain. A data item can be a large, complex object made up of many components, each relevant to some subset of the servers in the domain. Because these objects can be the units of consistency, it may be desirable to have a few large objects, rather than several tiny objects. As an example, a single data item or object can represent all configuration information for a system, including code files such as a config.xml file or an applicaiton-EAR file. A given component in the data item can, for example, be relevant to an individual server as to the number of threads, can be relevant to a cluster as to the deployed services, or can be relevant to the entire domain regarding security certificates. A delta between two versions can consist of new values for some or all of these components. For example, the components may include all enterprise Java beans deployed on members of the domain. A delta may include changes to only a subset of these Java beans.

[0044] The "scope" of a delta can refer to the set of all servers with a relevant component in the delta. An Admin server in accordance with the present invention may be able to interpret a configuration change in order to determine the scope of the delta. The DRS system on the master may need to know the scope in order to send the data to the appropriate slaves. It might be a waste of time and resources to send every configuration update to every server, when a master may only need to only touch a subset of servers in each update.

[0045] To control distribution, the master user can provide the scope of each update along with the delta between successive versions. A scope may be represented as a set of names, referring to servers and/or clusters, which may be taken from the same namespace within a domain. In one embodiment, the DRS uses a resolver module to map names to addresses. A cluster name can map to the set of addresses of all servers in that cluster. These addresses can be relative, such as to a virtual machine. The resolver can determine whether there is an intervening firewall, and return either an "inside" or "outside" address, relating to whether the server

is "inside the firewall" as is known and used in the art. An Admin server or other server can initialize the corresponding resolver with configuration data.

[0046] Along with the unique, long-lived domain identifier (DID) for each managed data item, each version of a data item can also have a long-lived version number. Each version number can be unique to an update attempt, such that a server will not improperly update or fail to update due to confusion as to the proper version. Similarly, the version number for an aborted two-phase distribution may not be re-used. The master may be able to produce a delta between two arbitrary versions given just the version numbers. If the master cannot produce such a delta, a complete copy of the data or application may be provided.

[0047] It may be desirable to keep the data replication service as generic as possible. A few assumptions may therefore be imposed upon the users of the system. The system may rely on, for example, three primary assumptions:

- the system may include a way to increment a version number
- the system may persistently store the version number on the master as well as the slave
- the system may include a way to compare version numbers and determine equality

These assumptions may be provided by a user-level implementation of a DRS interface, such as an interface "VersionNumber." Such an interface may allow a user to provide a specific notion and implementation of the version number abstraction, while ensuring that the system has access to the version number attributes. In Java, for example, a VersionNumber interface may be implemented as follows:

```
package weblogic.drs;  
public interface VersionNumber extends Serializable {  
    VersionNumber increment();  
    void persist() throws Exception;  
    boolean equals(VersionNumber anotherVN);  
    boolean strictlyGreaterThan(VersionNumber anotherVN);  
}
```


A simplistic implementation of this abstraction that a user could provide to the system would be a large, positive integer. The implementation may also ensure that the system can transmit delta information via the network from the master to the slaves, referred to in the art as being "serializable."

5 **[0048]** If using the abstraction above, it may be useful to abstract from a notion of the detailed content of a delta at the user level. The system may require no knowledge of the delta information structure, and in fact may not even be able to determine the structure. The
10 implementation of the delta can also be serializable, ensuring that the system can transmit delta version information via the network from the master to the slaves.

[0049] It may be desirable to have the master persistently store the copy of record for each data item, along with the appropriate DID and
15 version number. Before beginning a two-phase distribution, the master can persistently store the proposed new version number to ensure that it is not reused, in the event the master fails. A slave can persistently store the latest copy of each relevant data item along with its DID and version number. The slave can also be configured to do the necessary caching,
20 such that the slave may have to get the data or protocol every time. This may not be desirable in all cases, but may be allowed in order to handle certain situations that may arise.

[0050] A system in accordance with the present invention may further include concurrence restrictions. For instance, certain operations
25 may not be permitted during a two-phase distribution of an update for a given DID over a given scope. Such operations may include a one- or two-phase update, such as a modification of the membership of the scope on the same DID, over a scope with a non-empty intersection.

[0051] In at least one embodiment, the master DRS regularly
30 multicasts heartbeats, or packets of information, to the slave DRS on each server in the domain. For each DID, a heartbeat may contain a window of

information about the most recent update(s), including each update version number, the scope of the delta with respect to the previous version, and whether the update was committed or aborted. Information about the current version may always be included. Information about older versions
5 can also be used to minimize the amount of traffic back to the master, and not for correctness or liveness.

[0052] With the inclusion of older version information in a delta, the slave can commit that portion of the update it was expecting upon the prepare, and ask for a new delta to handle more recent updates.
10 Information about a given version can be included for at least some fixed, configurable number of heartbeats, although rapid-fire updates may cause the window to increase to an unacceptable size. In another embodiment, information about an older version can be discarded once a master determines that all slaves have received the update.

15 **[0053]** Multicast heartbeats may have several properties to be taken into consideration. These heartbeats can be asynchronous or "one-way". As a result, by the time a slave responds to a heartbeat, the master may have advanced to a new state. Further, not all slaves may respond at exactly the same time. As such, a master can assume that a slave has no
20 knowledge of its state, and can include that which the delta is intended to update. These heartbeats can also be unreliable, as a slave may miss arbitrary sequences of heartbeats. This can again lead to the inclusion of older version information in the heartbeats. In one embodiment, heartbeats are received by a slave in the order they were sent. For example, a slave
25 may not commit version seven until it has committed version six. The server may wait until it receives six, or it may simply throw out six and commit seven. This ordering may eliminate the possibility for confusion that might be created by versions going backwards.

[0054] As mentioned above, the domains may also utilize clustering,
30 as shown in Figure 3 (Properties of Multicast Heartbeats slide). The general network topology for this embodiment is a collection of multicast

islands, connected to a hub island containing the master. Multicast traffic may be forwarded point-to-point outward from the hub. Small deltas that may be distributed in the one-phase method may be directly transmitted over multicast. In all other cases, deltas may be transmitted over point-to-point links. A tree-structured, point-to-point forwarding scheme may be overlaid on the hub-and-spoke multicast structure to reduce the bottleneck at the master.

[0055] In the domain diagram **300** of **Figure 3**, one or more of the managed servers **302** can be grouped into a multicast island, also referred to as a cluster **304**. An Admin server **306** for the domain **308** acts as the master of the hub island **312**, and is the entry point to the domain, such as through a browser **310**. The Admin server **306** contacts one of the managed servers in the cluster, referred to as the cluster master. The Admin server in this embodiment can multicast a delta or message to each cluster master, with each cluster master then forwarding that delta or message by multicast to the other managed servers in that cluster. The cluster masters may not own any configuration information, instead receiving the information from the Admin server. In the event that a cluster master goes offline or crashes, another managed server in the domain may take over as cluster master. In this event, a mechanism may be put in place to prevent the offline server from coming back into the cluster as a second cluster master. This may be handled by the cluster or system infrastructure.

[0056] There can also be more than one domain. In this case, there can be nested domains or "syndicates." Information can be spread to the domain masters by touching each domain master directly, as each domain master can have the ability to push information to the other domain masters. It may, however, be undesirable to multicast to domain masters.

[0057] In one-phase distribution, a master user can make a downcall in order to trigger the distribution of an update. Such a downcall can take the form of:

startOnePhase(DID, newVerNum, scope)

where *DID* is the ID of the data item or object that was updated, *newVerNum* is the new version number of the object, and *scope* is the scope to which the update applies. The master DRS may respond by
5 advancing to the new version number, writing the new number to disk, and including the information in subsequent heartbeats.

[0058] When a slave DRS receives a heartbeat, it can determine whether it needs a pull by analyzing the window of information relating to recent updates of interest. If the slave's current version number is within
10 the window and the slave is not in the scope of any of the subsequent committed updates, it can simply advance to the latest version number without pulling any data. This process can include the trivial case where the slave is up-to-date. Otherwise, the slave DRS may make a point-to-point call for a delta from the master DRS, or another similar request, which
15 may take the form of:

createDelta(DID, curVerNum)

where *curVerNum* is the current number of the slave, which will be sent back to the domain master or cluster master. To handle this request, the master DRS may make an upcall, such as *createDelta(curVerNum)*. This
20 upcall may be made through the appropriate listener in order to obtain the delta and the new version number, and return them to the slave DRS. The new version number should be included, as it may have changed since the slave last received the heartbeat. The delta may only be up to the most recently committed update. Any ongoing two-phase updates may be
25 handled through a separate mechanism. The slave DRS may then make an upcall to the slave user, such as *commitOnePhase(newVerNum, delta)* and then advance to the new version number.

[0059] In order to trigger a two-phase update distribution, the master user can make a downcall, such as *startTwoPhase(DID, oldVerNum, newVerNum, delta, scope, timeout)*, where *DID* is the ID of the data item
30 or object to be updated, *oldVerNum* is the previous version number,

newVerNum is the new version number (one step from the previous version number), *delta* is the delta between the successive versions to be pushed, *scope* is the scope of the update, and *timeout* is the maximum time-to-live for the job. Because the "prepare" and "commit" are synchronous, it may be desirable to set a specific time limit for the job. The previous version number may be included to that a server on a different version number will not take the delta.

[0060] The master DRS in one embodiment goes through all servers in the scope and makes a point-to-point call to each slave DRS, such as *prepareTwoPhase(DID, oldVerNum, newVerNum, delta, timeout)*. The slave can then get the appropriate timeout value. Point-to-point protocol can be used where the delta is large, such as a delta that includes binary code. Smaller updates, which may for example include only minor configuration changes such as modifications of cache size, can be done using the one-phase method. This approach can be used because it may be more important that big changes like application additions get to the servers in a consistent fashion. The master can alternatively go to cluster masters, if they exist, and have the cluster masters make the call. Having the master proxy to the cluster masters can improve system scalability.

[0061] In one embodiment, each call to a slave or cluster master produces one of four responses, such as "Unreachable", "OutOfSync", "Nak", and "Ack", which are handled by the master DRS. If the response is "Unreachable", the server in question cannot be reached and may be queued for retry. If the response is "OutOfSync", the server may be queued for retry. In the meantime, the server will attempt to sync itself by using a pull from the master, so that it may receive the delta upon retry. If the response is "NoAck", or no acknowledgment, the job is aborted. This response may be given when the server cannot accept the job. If the response is "Ack", no action is taken.

[0062] In order to prepare the slaves, a master DRS can call a method such as *prepareTwoPhase*. Upon receiving a "prepare" request

from the master DRS, the slave DRS can first check whether its current version number equals the old version number to be updated. If not, the slave can return an "OutOfSync" response. The slave can then pull a delta from the master DRS as if it had just received a heartbeat. Eventually, the master DRS can retry the *prepareTwoPhase*. This approach may be more simple than having the master push the delta, but may require careful configuration of the master. The configuring of the master may be needed, as waiting too long for a response can cause the job to timeout. Further, not waiting long enough can lead to additional requests getting an "OutOfSync" response. It may be preferable to trigger the retry upon completion of the pull request from the slave.

[0063] If the slave is in sync, the slave can make an upcall to the client layer on the slave side, as deep into the server as possible, such as *prepareTwoPhase(newVerNum, delta)*. The resulting "Ack" or "Nak" that is returned can then be sent to the master DRS. If the response was an "Ack", the slave can go into a special prepared state. If the response was a "Nak", the slave can flush any record of the update. If it were to be later committed for some reason, the slave can obtain it as a one-phase distribution, which may then fail.

[0064] If the master DRS manages to collect an "Ack" from every server within the timeout period, it can make a commit upcall, such as *twoPhaseSucceeded(newVerNum)*, and advance to the new version number. If the master DRS receives a "Nak" from any server, or if the timeout period expires, the master DRS can make an abort upcall, such as *twoPhaseFailed(newVerNum, reason)*, and leave the version number unchanged. Here, *reason* is an exception, containing a roll-up of any "Nak" responses. In both cases, the abort/commit information can be included in subsequent heartbeats.

[0065] At any time, the master DRS can make a cancel downcall, such as *cancelTwoPhase(newVerNum)*. The master DRS can then handle

this call by throwing an exception, if the job is not in progress, or acting as if an abort is to occur.

- 5 **[0066]** If a prepared slave DRS gets a heartbeat indicating the new version was committed, the slave DRS can make an upcall, such as *commitTwoPhase(newVerNum)*, and advance to the new version number. If a prepared slave DRS instead gets a heartbeat indicating the new version was aborted, the slave can abort the job. The slave can also abort the job when the slave gets a heartbeat where the window has advanced beyond the new version, the slave gets a new *prepareTwoPhase* call on the same data item, or the slave times out the job. In such a case, the slave can make an upcall, such as *abortTwoPhase(newVerNum)*, and leave the version number unchanged. This is one way to ensure the proper handling of situations such as where a master server fails after the slaves were prepared but before the slaves commit.
- 10
- 15 **[0067]** The foregoing description of preferred embodiments of the present invention has been provided for the purposes of illustration and description. It is not intended to be exhaustive or to limit the invention to the precise forms disclosed. Obviously, many modifications and variations will be apparent to the practitioner skilled in the art. The embodiments were chosen and described in order to best explain the principles of the invention and its practical application, thereby enabling others skilled in the art to understand the invention for various embodiments and with various modifications that are suited to the particular use contemplated. It is intended that the scope of the invention be defined by the following claims and their equivalence.
- 20

CLAIMS

What is claimed is:

- 1 1. A method for replicating data from a master server to a slave server over
2 a network, the method comprising the steps of:
3 sending a packet of information from the master server to the
4 slave server, the information relating to a change in
5 the data stored on the master server and containing a
6 version number for the present state of the data;
7 allowing the slave server to determine whether the data on
8 the slave server has been updated to correspond to
9 the version number contained in the packet; and
10 requesting a delta be sent from the master server to the
11 slave server if the data on the slave server does not
12 correspond to the version number contained in the
13 packet, the delta containing information needed to
14 update the slave server.
- 1 2. A method according to claim 1, further comprising:
2 storing an original copy of the data on the master server.
- 1 3. A method according to claim 1, further comprising:
2 persistently caching the data on a local disk for each slave server.
- 1 4. A method according to claim 1, further comprising:
2 determining a unique version number for the current state of the
3 data on the master server if the data has changed.
- 1 5. A method for replicating data from a master server to a slave server over
2 a network, the method comprising the steps of:

3 sending a version number from the master server to the
4 slave server, the version number relating to the
5 present state of the data stored on the master server;
6 allowing the slave server to determine whether the slave
7 server has been updated to reflect the present state of
8 the data corresponding to the version number sent
9 from the master server; and
10 requesting a delta be sent from the master server to the
11 slave server if the slave server does not correspond to
12 the version number sent by the master, the delta
13 containing information needed to update the slave
14 server.

1 6. A method according to claim 5, further comprising:

2 sending the delta from the master server to the slave server.

1 7. A method according to claim 5, further comprising:

2 committing the delta to the slave server.

1 8. A method according to claim 5, further comprising:

2 updating the version number of the slave server after committing the
3 delta.

1 9. A method according to claim 5, further comprising:

2 periodically sending the version number from the master server to
3 a slave server.

1 10. A method according to claim 5, further comprising:

2 sending the version number to a slave server until the slave server
3 acknowledges receipt of the version number.

1 11. A method according to claim 5, further comprising:
2 including data with the version number that is necessary to update
3 a slave server.

1 12. A method according to claim 11, further comprising:
2 committing the data necessary to update the slave server as soon
3 as it is received.

1 13. A method according to claim 5, further comprising:
2 determining the scope of the delta before sending it from the master
3 server.

1 14. A method for replicating data over a network including a master
2 server and at least one slave server, the method comprising the
3 steps of:
4 sending a packet of information from a master server to each
5 slave server on the network, the information relating to
6 a change in the data stored on the master server and
7 containing a current version number for the present
8 state of the data, the information further relating to
9 previous changes in the data and a version number
10 for each previous change;
11 allowing each slave server to determine whether the slave
12 server has been updated to correspond to the current
13 version number;
14 allowing each slave server to commit the information if the
15 slave server has not missed a previous change; and
16 allowing each slave server having missed a previous change
17 to request that previous change be sent from the
18 master server to the slave server before the slave
19 server commits the packet of information.

- 1 15. A method according to claim 14, further comprising:
2 committing the packet of information to a slave server.
- 1 16. A method according to claim 14, further comprising:
2 aborting the commit of the packet of information if a slave server
3 cannot commit the update.
- 1 17. A method according to claim 14, further comprising:
2 determining the scope of the delta before sending it from the master
3 server.
- 1 18. A method according to claim 14, further comprising:
2 including the scope of each the previous changes in the delta.
- 1 19. A method for replicating data over a network including a master server
2 and at least one slave server, the method comprising the steps of:
3 sending a packet of information from a master server to each
4 slave server on the network, the information relating to
5 a change in the data stored on the master server and
6 containing a prior version number for the prior state
7 and a new version number for the new state of the
8 data, the information further relating to previous
9 changes in the data and a previous version number
10 for each previous change;
11 allowing each slave server to determine whether the data on
12 the slave server corresponds to the prior version
13 number contained in the packet;
14 allowing each slave server to commit the packet of
15 information if the data on the slave server corresponds
16 to the prior version number contained in the packet,

17 the commit also updating the version of the slave
18 server to the new version number; and
19 allowing each slave server not corresponding to the prior
20 version number to request that a delta be sent from
21 the master server containing the information
22 necessary to update the slave to the prior version
23 number before the slave server commits the packet of
24 information.

1 20. A method for replicating data over a network including a master server
2 and at least one slave server, the method comprising the steps of:
3 sending a packet of information from a master server to each
4 slave server on the network, the information relating to
5 a change in the data stored on the master server and
6 containing a version number for the prior state and a
7 version number for the new state of the data, the
8 information further relating to previous changes in the
9 data and a version number for each previous change;
10 allowing each slave server to determine whether the data on
11 the slave server corresponds to the prior version
12 number contained in the packet;
13 allowing each slave server to commit the packet of
14 information if the data on the slave server corresponds
15 to the prior version number contained in the packet,
16 the commit also updating the version of the slave
17 server to the new version number; and
18 allowing each slave server not corresponding to the prior
19 version number to request that a delta be sent from
20 the master server containing the information
21 necessary to update the slave to the new version
22 number.

1 21. A method for replicating data from a master server to at least one slave
2 server over a network, the method comprising the steps of:
3 sending a packet of information from the master server to a
4 slave server, the information relating to a change in
5 the data stored on the master server and containing a
6 version number for the present state of the data;
7 receiving the packet of information to a slave server;
8 allowing the slave server to determine whether the slave
9 server has been updated to correspond to the version
10 number contained in the packet, and to further
11 determine whether the slave server can process the
12 packet of information if needed to update to
13 correspond to the version number contained in the
14 packet;
15 sending a signal from the slave server to the master server,
16 the signal indicating whether the slave server needs to
17 be updated and whether the slave server can process
18 the update; and
19 sending a response signal from the master server to the
20 slave server indicating whether the slave server
21 should commit to the information contained in the
22 packet; and
23 committing the packet of information to the slave server if so
24 indicated by the response signal.

1 22. A method according to claim 21, further comprising:
2 determining whether each of the at least one slave server can
3 commit the data.

1 23. A method according to claim 21, further comprising:

2 determining whether each of the at least one slave server has sent
3 a response back to the master server.

1 24. A method according to claim 21, further comprising:
2 determining whether any of the at least one slave server can commit
3 the data.

1 25. A method according to claim 21, further comprising:
2 committing the data only if each of the at least one slave server can
3 process the commit.

1 26. A method according to claim 21, further comprising:
2 aborting the data only if any of the at least one slave server cannot
3 process the commit.

1 27. A method according to claim 21, further comprising:
2 committing the data to those slaves that are able to process the
3 commit.

1 28. A method according to claim 21, further comprising:
2 multicasting the update to any of the at least one slave server that
3 were not able to process the commit.

1 29. A method according to claim 21, further comprising:
2 heartbeating the new version number to any of the at least one
3 slave server that were not able to process the commit.

1 30. A method according to claim 21, further comprising:
2 requesting a delta be sent to a slave server that was not able to
3 process the commit.

- 1 31. A method for replicating data over a network, the method comprising
2 the steps of:
- 3 (a) determining whether the replication should be accomplished in
4 a one or two phase method;
- 5 (b) sending replication information determined to be accomplished
6 in a one phase method by:
7 sending a packet of information from the master server to the
8 slave server, the information relating to a change in
9 the data stored on the master server and containing a
10 version number for the present state of the data;
11 receiving the packet of information to a slave server;
12 allowing the slave server to determine whether the data on
13 the slave server has been updated to correspond to
14 the version number; and
15 requesting a delta be sent from the master server to the
16 slave server if the slave server does not correspond to
17 the version number, the delta containing information
18 needed to update the slave server;
- 19 (c) sending replication information determined to be accomplished
20 in a two phase method by:
21 sending a packet of information from the master server to the
22 slave server, the information relating to a change in
23 the data stored on the master server and containing a
24 version number for the present state of the data;
25 allowing the slave server to determine whether the slave
26 server has been updated to correspond to the version
27 number, and to further determine whether the slave
28 server can process the packet of information;
29 sending a signal from the slave server to the master server
30 indicating whether the slave server needs to be

31 updated and whether the slave server can process the
32 packet of information;
33 sending a response signal from the master server to the
34 slave server indicating whether the slave server
35 should commit to the packet of information; and
36 committing the packet of information to the slave server if so
37 indicated by the response signal.

1 32. A method for replicating data over a network, the method comprising
2 the steps of:
3 (a) determining whether replication should be accomplished in a one
4 or two phase method;
5 (b) sending data to be replicated in a one phase method by:
6 sending a version number for the current state of the data
7 from a master server to a slave server;
8 requesting a delta be sent from the master server to the
9 slave server if the data on the slave server does not
10 correspond to the version number; and
11 (c) sending data to be replicated in a two phase method by:
12 sending a packet of information from the master server to a
13 slave server;
14 determining whether the slave server can process the packet
15 of information; and
16 committing the packet of information to the slave server if the
17 slave server can process the packet of information.

1 33. A method for replicating data from a master to a plurality of slaves on
2 a network, the method comprising the steps of:
3 (a) determining whether replication should be accomplished in a one
4 or two phase method;
5 (b) sending data to be replicated in a one phase method by:

6 sending a version number for the current state of the data
7 from the master to each slave; and
8 requesting a delta be sent from the master to each slave
9 containing data that does not correspond to the
10 version number;
11 (c) sending data to be replicated in a two phase method by:
12 sending a packet of information from the master to each
13 slave; and
14 committing the packet of information to the slaves if each of
15 the plurality of slaves can process the packet of
16 information.

1 34. A method for replicating data from a master to a plurality of slaves on
2 a network using one and two phase methods, the method comprising the
3 steps of:

4 (a) sending data to be replicated in a one phase method by sending
5 a version number for the current state of the data from the
6 master to each slave so that each slave may request a delta
7 to be sent from the master to the slave to update the data on
8 the slave; and
9 (b) sending data to be replicated in a two phase method by sending
10 a packet of information from the master to each slave, the
11 packet of information to be committed by each slave if every
12 slave is able to commit the packet of information.

1 35. A method for replicating data on a clustered network using one and two
2 phase methods, each network cluster containing a cluster master and at
3 least one cluster slave, the method comprising the steps of:

4 (a) sending data to be replicated in a one phase method by sending
5 a version number for the current state of the data from a first

6 cluster master to all other cluster masters so the other cluster
7 masters may each request a delta; and

8 (b) sending data to be replicated in a two phase method by sending
9 a packet of information from the first cluster master to each
10 other cluster master, the packet of information to be
11 committed by the other cluster masters if the other cluster
12 masters are able to commit the packet of information.

1 36. A method according to claim 35, further comprising:

2 sending the data from each cluster master to each cluster slave in
3 the cluster with that cluster master by a one-phase method.

1 37. A method according to claim 10, further comprising:

2 sending the data from each cluster master to each cluster slave in
3 the cluster with that cluster master by a two-phase method.

1 38. A computer-readable medium, comprising:

2 (a) means for sending a packet of information from a master server
3 to each slave server on the network, the information relating
4 to a change in the data stored on the master server and
5 containing a current version number for the present state of
6 the data, the information further relating to previous changes
7 in the data and a version number for each previous change;

8 (b) means for allowing each slave server to determine whether the
9 slave server has been updated to correspond to the current
10 version number;

11 (c) means for allowing each slave server to commit the information
12 if the slave server has not missed a previous change; and

13 (d) means for allowing each slave server having missed a previous
14 change to request that previous change be sent from the

15 master server to the slave server before the slave server
16 commits the packet of information.

1 39. A computer program product for execution by a server computer for
2 replicating data over a network, comprising:

3 (a) computer code for sending a packet of information from a master
4 server to each slave server on the network, the information
5 relating to a change in the data stored on the master server
6 and containing a current version number for the present state
7 of the data, the information further relating to previous
8 changes in the data and a version number for each previous
9 change;

10 (b) computer code for allowing each slave server to determine
11 whether the slave server has been updated to correspond to
12 the current version number;

13 (c) computer code for allowing each slave server to commit the
14 information if the slave server has not missed a previous
15 change; and

16 (d) computer code for allowing each slave server having missed a
17 previous change to request that previous change be sent
18 from the master server to the slave server before the slave
19 server commits the packet of information.

1 40. A system for replicating data over a network, comprising:

2 (a) means for sending a packet of information from a master
3 server to each slave server on the network, the
4 information relating to a change in the data stored on
5 the master server and containing a current version
6 number for the present state of the data, the
7 information further relating to previous changes in the
8 data and a version number for each previous change;

(b) means for allowing each slave server to determine whether the slave server has been updated to correspond to the current version number;

(c) means for allowing each slave server to commit the information if the slave server has not missed a previous change; and

(d) means for allowing each slave server having missed a previous change to request that previous change be sent from the master server to the slave server before the slave server commits the packet of information.

41. A computer system comprising:

a processor;

object code executed by said processor, said object code configured to:

(a) send a packet of information from a master server to each slave server on the network, the information relating to a change in the data stored on the master server and containing a current version number for the present state of the data, the information further relating to previous changes in the data and a version number for each previous change;

(b) allow each slave server to determine whether the slave server has been updated to correspond to the current version number;

(c) allow each slave server to commit the information if the slave server has not missed a previous change; and

(d) allow each slave server having missed a previous change to request that previous change be sent from the master server to the slave server before the slave server commits the packet of information.

21

22

42. A system for replicating data over a network, comprising:

23

a. a master server containing an original copy of the data, said master server comprising:

24

25

i. a master user layer adapted to start a data replication process by calling a start method, the master user layer further adapted to send information relating to the original copy of the data;

26

27

28

29

ii. a master service layer containing the start method and adapted to receive the call from the master user layer and the information relating to the original copy of the data, the master service layer further adapted to create and send a data replication packet containing at least some of the information relating to the original copy of the data;

30

31

32

33

34

35

36

b. a slave server adapted to store a copy of the data from the master server, the slave server comprising:

37

38

i. a slave service layer adapted to receive the data replication packet from the master service layer and process the data replication packet, the slave service layer further adapted to send information relating to the data replication packet; and

39

40

41

42

43

ii. a slave user layer adapted to receive the information relating to the data replication packet from the slave service layer, the slave user layer adapted to store the information in the data replication packet.

44

45

46

1

43. A system according to claim 42, wherein said master user layer is in communication with at least one of a master user and a master user device.

2

3

1 44. A system according to claim 42, wherein said master user layer is
2 adapted to send information relating to the original copy of the data
3 in the form of a delta, the delta containing information relating to
4 changes between a previous state and the current state of the
5 original copy of the data.

1 45. A system according to claim 42, wherein said master user layer is
2 adapted to update the original copy of the data.

1 46. A system according to claim 42, wherein said master user layer is
2 adapted to send a roll-back message indicating that a change to the
3 original copy of the data should not be replicated on a slave server.

1 47. A system according to claim 42, wherein said master user layer is
2 adapted to set a timeout value for the replication.

1 48. A system according to claim 42, wherein said master user layer is
2 adapted to create a delta between the present state of the original
3 copy of the data and the prior state of the original copy of the data.

1 49. A system according to claim 42, wherein said master user layer is
2 adapted to create a delta between the present state of the original
3 copy of the data and a previous state of the original copy of the
4 data.

1 50. A system according to claim 42, wherein said master user layer is
2 adapted to generate a unique version number for each state of the
3 original copy of the data.

1 51. A system according to claim 42, wherein said master service layer
2 is adapted to multicast the data replication packet.

- 1 52. A system according to claim 42, wherein said master service layer
2 is adapted to heartbeat the data replication packet.
- 1 53. A system according to claim 42, wherein said master service layer
2 is adapted to include a version number in the data replication
3 packet.
- 1 54. A system according to claim 42, wherein said master service layer
2 is adapted to include information necessary to update the copy of
3 the data on the slave server to the current state of the original copy
4 of the data.
- 1 55. A system according to claim 42, wherein said master service layer
2 is further adapted to create and send a data replication packet
3 comprising a delta.
- 1 56. A system according to claim 42, wherein said master service layer
2 is further adapted to create and send a data replication packet
3 comprising a delta between successive states of the original copy of
4 the data.
- 1 57. A system according to claim 42, wherein said master service layer
2 is further adapted to create and send a data replication packet
3 comprising a delta between arbitrary states of the original copy of
4 the data.
- 1 58. A system according to claim 42, wherein said master service layer
2 is adapted to request a delta from the master user layer.

- 1 59. A system according to claim 42, wherein said master service layer
2 is adapted to send a commit message to a slave service layer.
- 1 60. A system according to claim 42, wherein said master service layer
2 is adapted to heartbeat a commit message to a slave service layer.
- 1 61. A system according to claim 42, wherein said master service layer
2 is adapted to multicast a commit message to a slave service layer.
- 1 62. A system according to claim 42, wherein said master service layer
2 is adapted to send an abort message to a slave service layer.
- 1 63. A system according to claim 42, wherein said master service layer
2 is adapted to heartbeat an abort message to a slave service layer.
- 1 64. A system according to claim 42, wherein said master service layer
2 is adapted to multicast an abort message to a slave service layer.
- 1 65. A system according to claim 42, wherein said slave user layer is in
2 communication with at least one of a slave user and a slave user
3 device.
- 1 66. A system according to claim 42, wherein said slave user layer is
2 adapted to check the current version number of data stored on the
3 slave server.
- 1 67. A system according to claim 42, wherein said slave user layer is
2 adapted to commit information relating to the data replication packet
3 to the data stored on the slave server.

- 1 68. A system according to claim 42, wherein said slave user layer is
2 adapted to abort an update to the data stored on the slave server.
- 1 69. A system according to claim 42, wherein said slave user layer is
2 adapted to process a prepare request contained in the data
3 replication packet.
- 1 70. A system according to claim 42, wherein said slave user layer is
2 adapted to send a response to the slave service layer relating to a
3 prepare request contained in the data replication packet.
- 1 71. A system according to claim 42, wherein said slave user layer is
2 adapted to persistently cache data on a local disk.
- 1 72. A system according to claim 42, wherein said slave user layer is
2 adapted to update the version number of the copy of the data on the
3 slave server.
- 1 73. A system according to claim 42, wherein said slave service layer is
2 adapted to request a delta from the master service layer.
- 1 74. A system according to claim 42, wherein said slave service layer is
2 adapted to request the current version number of the data stored on
3 the slave server from the slave user layer.
- 1 75. A system according to claim 42, wherein said slave service layer is
2 adapted to send a commit message to the slave user layer.
- 1 76. A system according to claim 42, wherein said slave service layer is
2 adapted to send an abort message to the slave user layer.

1 77. A method for replicating data from a master server to a slave server,
2 comprising:

3 78. sending a start call from a master user level to a master service level
4 on a master server, the start call containing information relating to
5 the current state of master data on the master server;

6 a. sending the information to a slave service layer on a slave
7 server, the slave service layer adapted to check a slave user
8 layer on the slave server to determine whether slave data on
9 the slave server has the current state;

10 b. sending a request for a delta from the slave service layer to
11 the master service layer, the master service layer adapted to
12 request and receive a delta from the master user layer;

13 c. sending a delta from the master service layer to the slave
14 service layer, the delta containing the information necessary
15 to bring the slave data up to the current state, the slave
16 service layer adapted to process the delta and send the
17 information to the slave user layer; and

18 d. updating the slave data using the slave user layer.

1 79. A method according to claim 77, further comprising:
2 determining a version number for the current state of the data using
3 the master user layer.

1 80. A method according to claim 77, further comprising:
2 sending the information to the slave service layer by multicasting.

1 81. A method according to claim 77, further comprising:
2 sending information to the slave service layer, the information
3 comprising a version number for the current state of the master data.

- 1 82. A method for replicating data from a master server to a slave server,
2 comprising:
3 a. sending a new delta from a master user level to a master
4 service level on a master server, the delta containing
5 information relating to a change from the prior state to the
6 current state in master data stored on the master server;
7 b. sending the new delta from the master service layer to a
8 slave service layer on a slave server, the slave service layer
9 adapted to check a slave user layer on the slave server to
10 determine whether the slave data on the slave server has the
11 current state;
12 c. sending a request for a syncing delta from the slave service
13 layer to the master service layer, the master service layer
14 adapted to request and receive a syncing delta from the
15 master user layer, the syncing delta containing information
16 necessary to update the slave data to the prior state of the
17 master data;
18 d. sending the syncing delta from the master service layer to the
19 slave service layer, the slave service layer adapted to
20 process the delta and send the information to the slave user
21 layer to be committed to the slave data; and
22 e. committing the information in the new delta to the slave data
23 using the slave user layer.

- 1 83. A method for replicating data from a master server to a slave server
2 over a network, the method comprising the steps of:
3 a. sending a version number from a master service layer to a
4 slave service layer relating to the present state of the original
5 copy of the data on the master server;

- 6 b. allowing a slave user layer to determine whether the data on
7 the slave server has been updated to correspond to the
8 version number; and
9 c. requesting a delta be sent from the master service layer to
10 the slave service layer if the data on the slave server does
11 not correspond to the version number.

1 84. A method according to claim 77, further comprising:
2 allowing the slave user layer to persistently cache the data on a local
3 disk for each slave server.

1 85. A method according to claim 77, further comprising:
2 allowing the master user layer to determine a unique version number
3 for the current state of the data on the master server.

1 86. A method according to claim 77, further comprising:
2 including data with the version number that is necessary for a slave
3 user layer to update the data on a slave server.

1 87. A method according to claim 77, further comprising:
2 committing the data necessary to update the slave server as soon
3 as it is received by the slave user layer.

1 88. A method for replicating data over a network including a master
2 server and at least one slave server, the method comprising the
3 steps of:
4 a. sending a packet of information from a master service layer
5 to a slave service layer on each slave server on the network,
6 the information relating to a change in the data stored on the
7 master server and containing a prior version number for the
8 prior state and a new version number for the new state of the

9 data, the information further relating to previous changes in
10 the data and a previous version number for each previous
11 change;
12 b. allowing a slave user layer on each slave server to determine
13 whether the data on the slave server corresponds to the prior
14 version number contained in the packet;
15 c. allowing each slave user layer to commit the packet of
16 information if the data on the slave server corresponds to the
17 prior version number contained in the packet, the commit
18 also updating the version of the slave server to the new
19 version number; and
20 d. allowing each slave user layer not corresponding to the prior
21 version number to request that a delta be sent from the
22 master service layer to the slave service layer corresponding
23 to that slave user layer, the delta containing the information
24 necessary to update the slave to the prior version number
25 before the slave service layer commits the packet of
26 information.

1 89. A method for replicating data from a master server to at least one
2 slave server over a network, the method comprising the steps of:
3 a. sending a packet of information from a master service layer
4 on the master server to the user service layer on a slave
5 server, the information relating to a change in the data stored
6 on the master server and containing a version number for the
7 present state of the data;
8 b. allowing the slave user layer on the server to determine
9 whether the slave server has been updated to correspond to
0 the version number contained in the packet, and to further
1 determine whether the slave user layer can process the

- 12 packet of information if needed to update to correspond to
13 the version number contained in the packet;
- 14 c. sending a signal from the slave service layer to the master
15 service layer, the signal indicating whether the slave server
16 needs to be updated and whether the slave server can
17 process the update;
- 18 d. sending a response signal from the master service layer to
19 the slave service layer indicating whether the slave user layer
20 should commit to the information contained in the packet;
21 and
- 22 e. committing the packet of information to the slave server if so
23 indicated by the response signal.

- 1 90. A computer-readable medium, comprising:
- 2 a. means for sending a version number from a master service layer
3 to a slave service layer relating to the present state of the
4 original copy of the data on the master server;
- 5 b. means for allowing a slave user layer to determine whether the
6 data on the slave server has been updated to correspond to
7 the version number; and
- 8 c. means for requesting a delta be sent from the master service
9 layer to the slave service layer if the data on the slave server
10 does not correspond to the version number.

- 1 91. A computer program product for execution by a server computer for
2 replicating data from a master server to a slave server over a
3 network, comprising:
- 4 a. computer code for sending a version number from a master
5 service layer to a slave service layer relating to the present
6 state of the original copy of the data on the master server;

- 7 b. computer code for allowing a slave user layer to determine
- 8 whether the data on the slave server has been updated to
- 9 correspond to the version number; and
- 10 c. computer code for requesting a delta be sent from the master
- 11 service layer to the slave service layer if the data on the slave
- 12 server does not correspond to the version number.

- 1 92. A system for replicating data over a network, comprising:
- 2 a. means for sending a version number from a master service layer
- 3 to a slave service layer relating to the present state of the
- 4 original copy of the data on the master server;
- 5 b. means for allowing a slave user layer to determine whether the
- 6 data on the slave server has been updated to correspond to
- 7 the version number; and
- 8 c. means for requesting a delta be sent from the master service
- 9 layer to the slave service layer if the data on the slave server
- 10 does not correspond to the version number.

- 1 93. A computer system comprising:
- 2 a processor;
- 3 object code executed by said processor, said object code configured
- 4 to:
- 5 a. send a version number from a master service layer to a
- 6 slave service layer relating to the present state of the
- 7 original copy of the data on the master server;
- 8 b. allow a slave user layer to determine whether the data on
- 9 the slave server has been updated to correspond to
- 10 the version number; and
- 11 c. request a delta be sent from the master service layer to
- 12 the slave service layer if the data on the slave server
- 13 does not correspond to the version number.

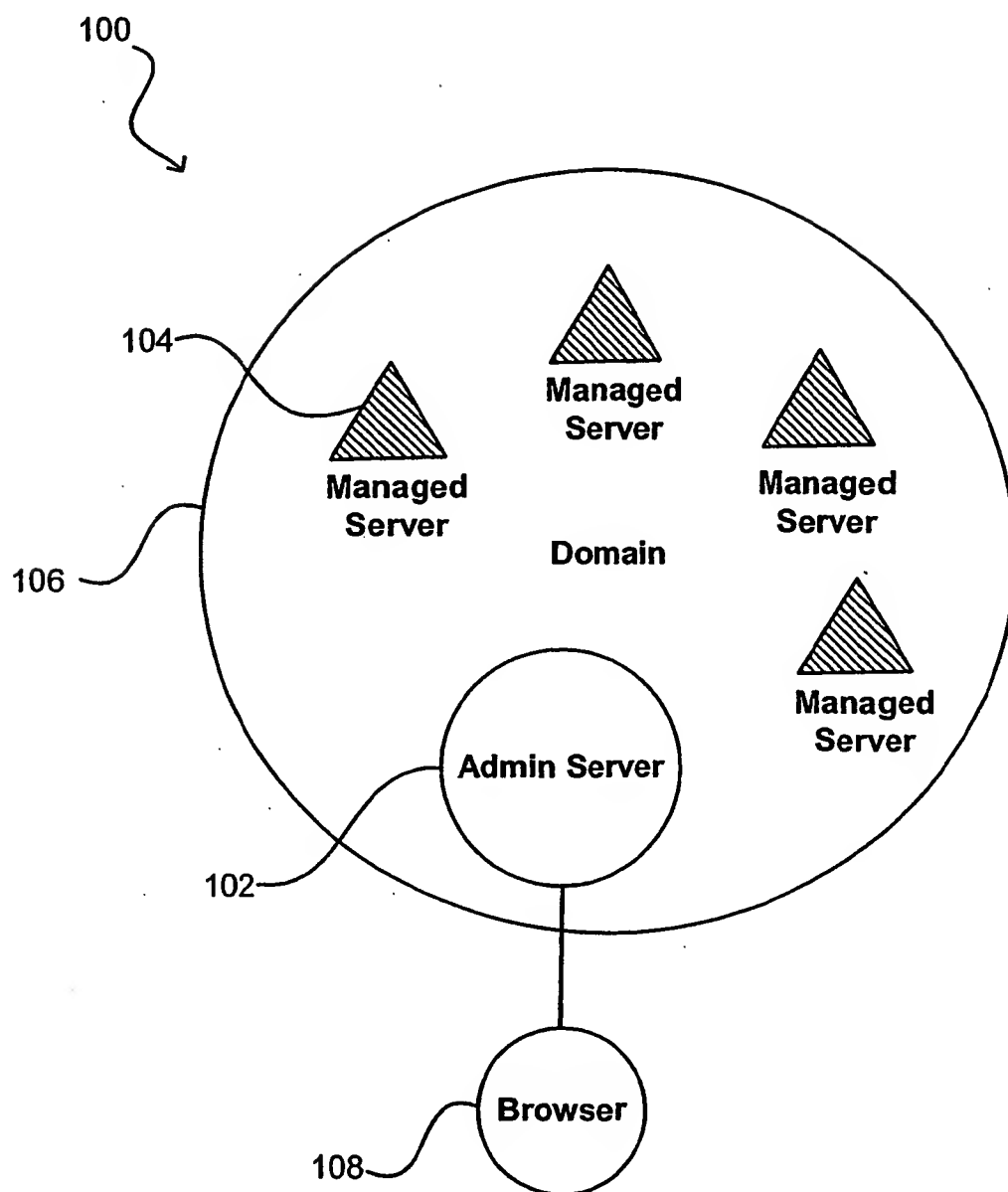


Figure 1

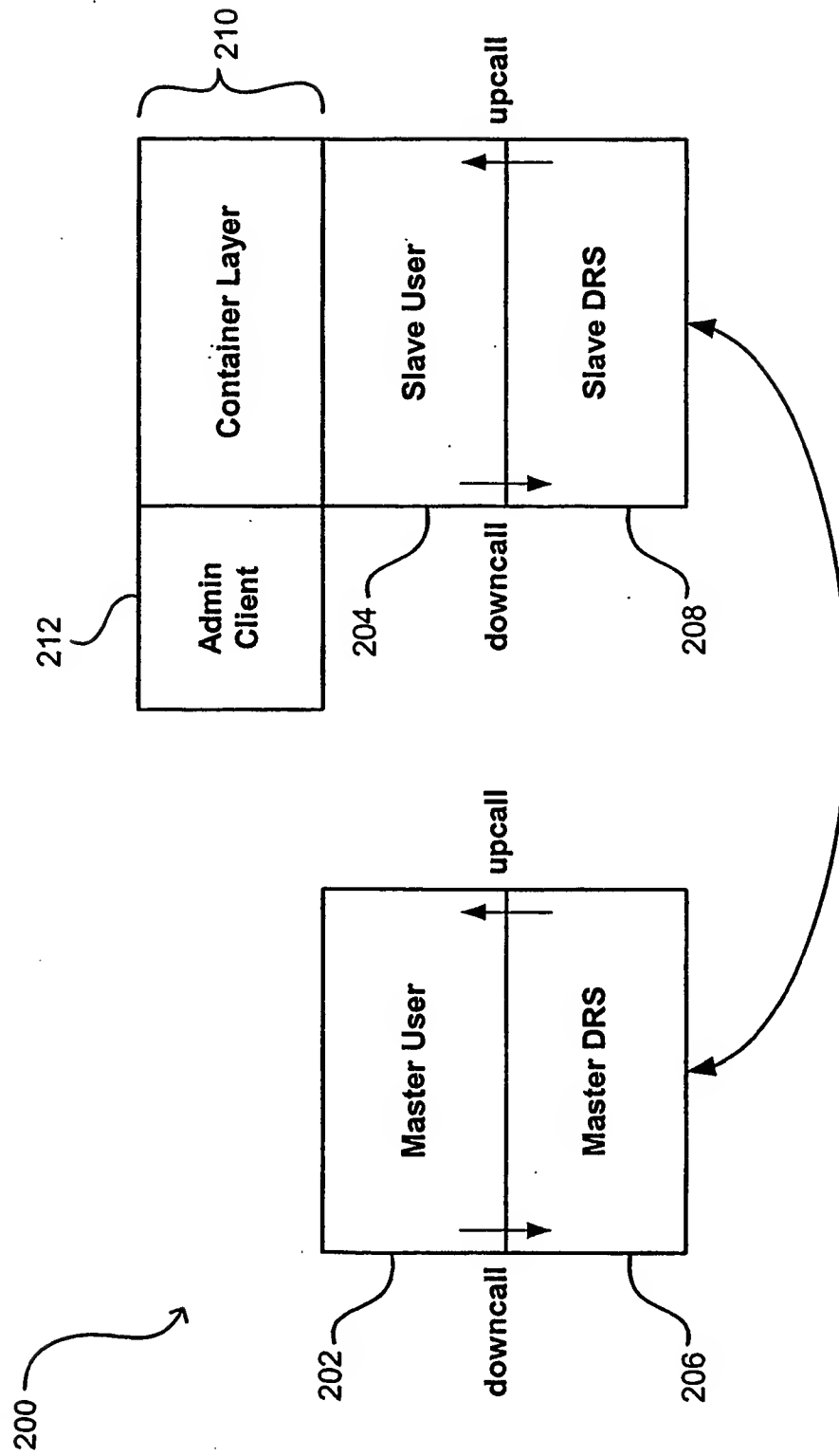


Figure 2

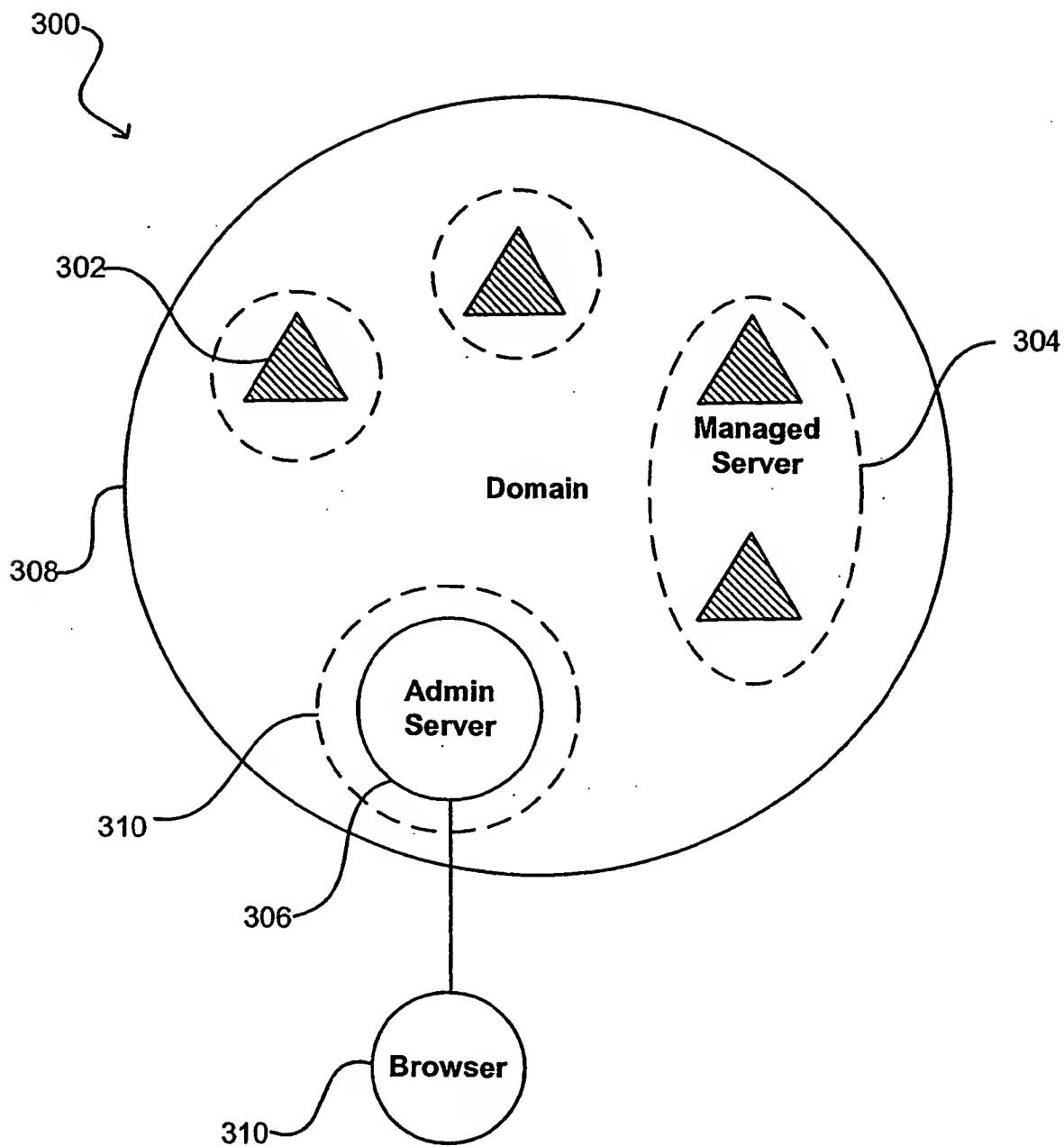


Figure 3

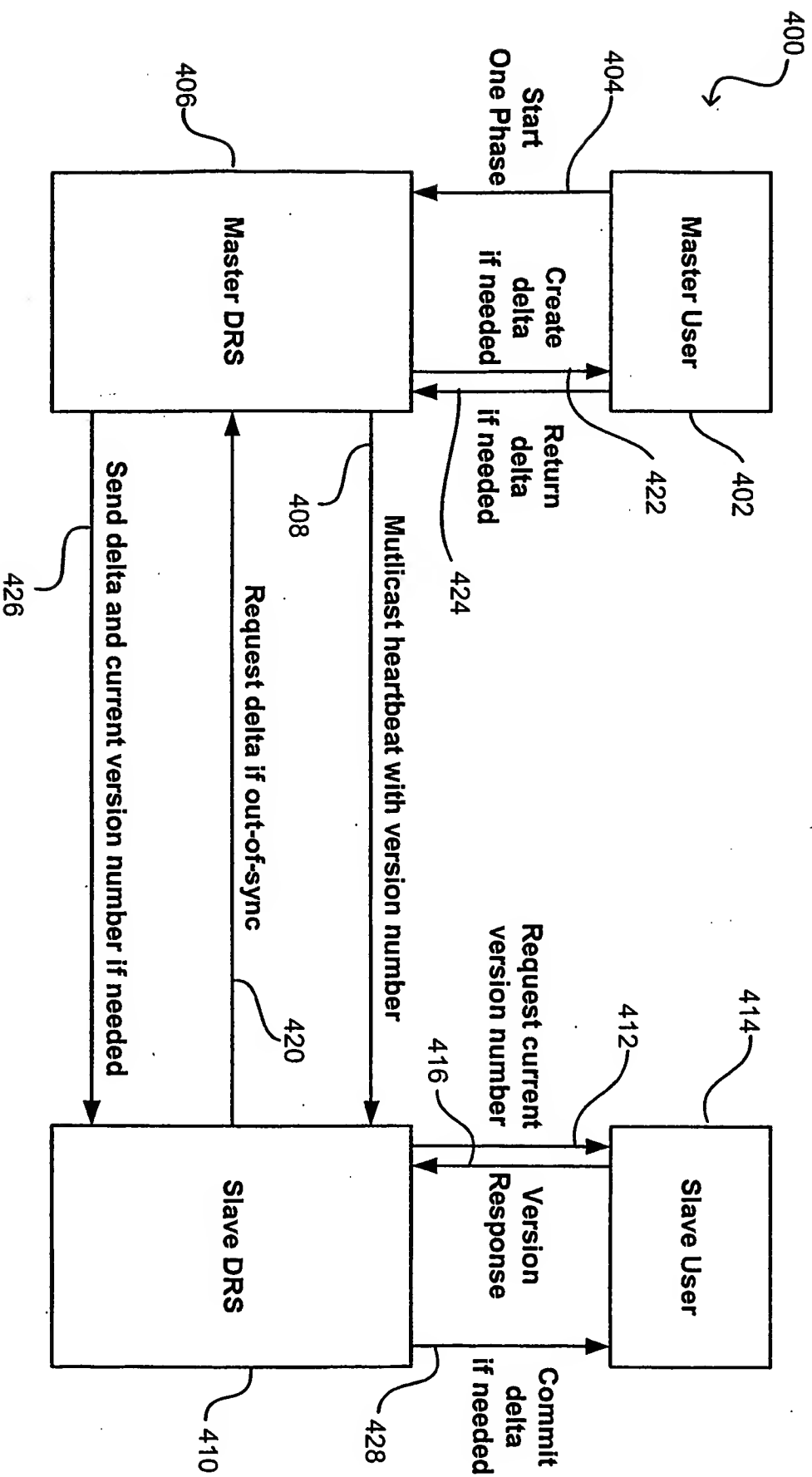


Figure 4

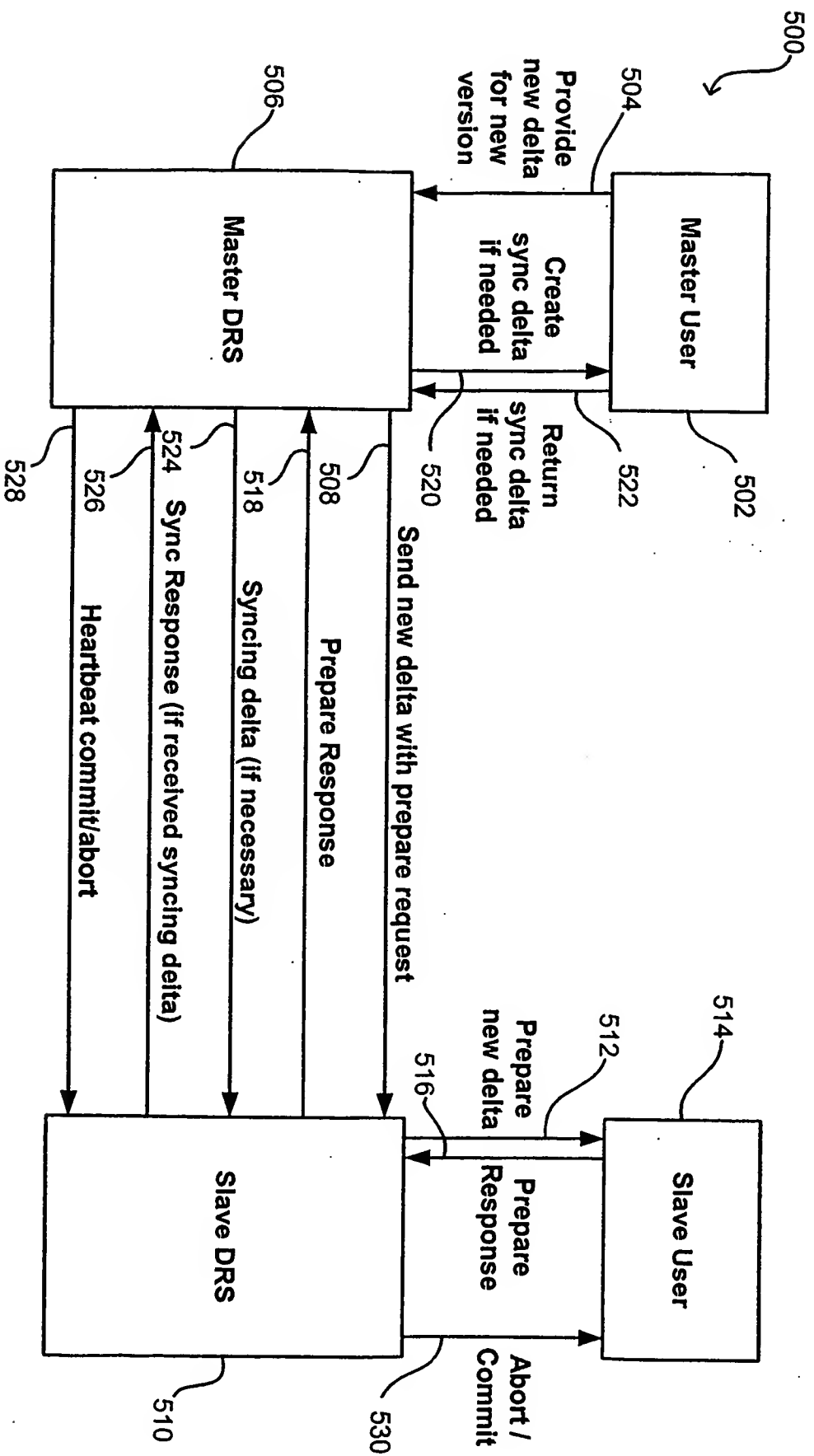
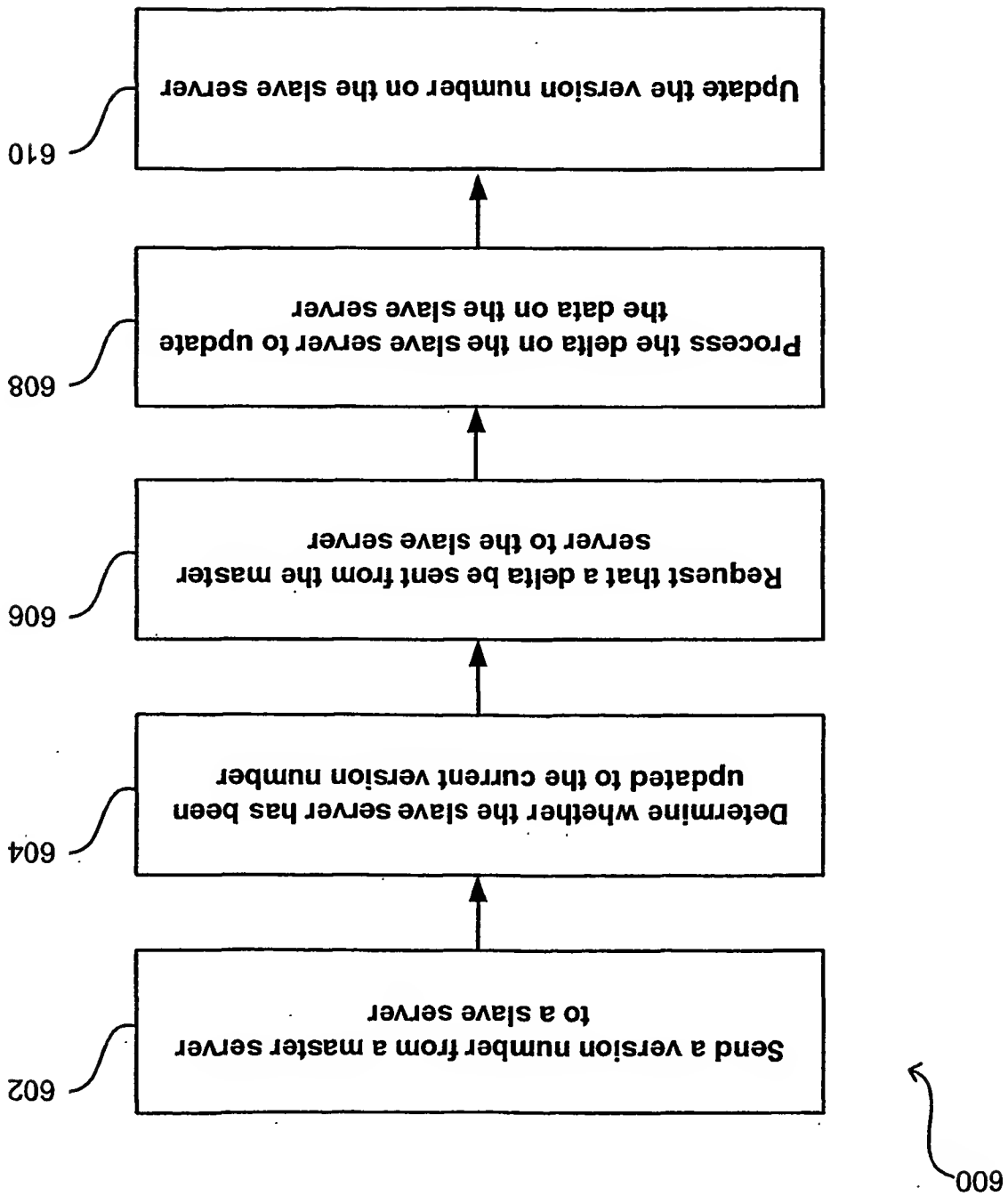


Figure 5

*Figure 6*

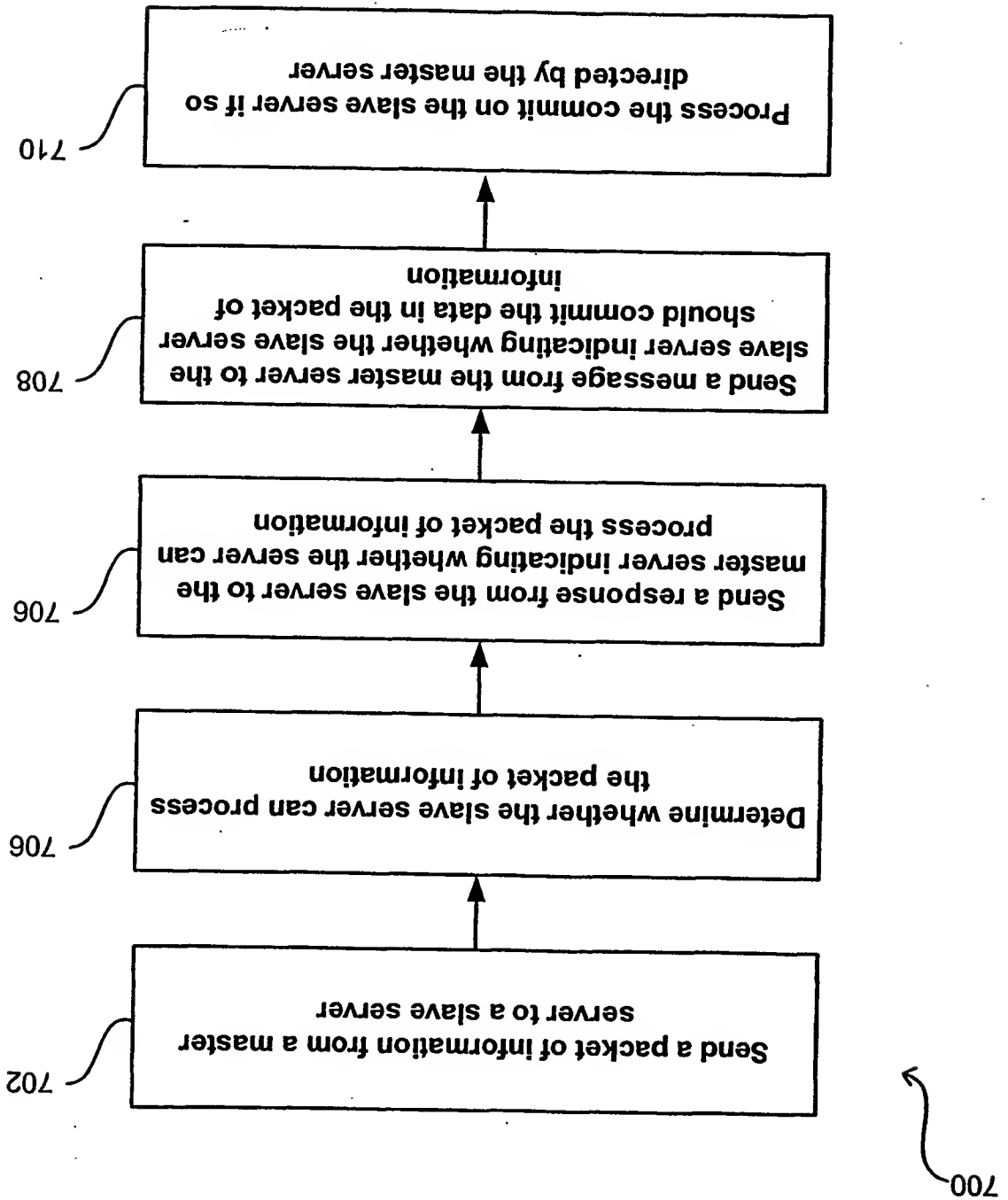


Figure 7

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
30 January 2003 (30.01.2003)

PCT

(10) International Publication Number
WO 03/009092 A3

(51) International Patent Classification⁷: **G06F 15/16**,
7/00, 17/30

SRINIVASAN, Ananthan, Bala; 1610 Sanchez Street,
San Francisco, CA 94131 (US).

(21) International Application Number: PCT/US02/22366

(74) Agents: **MEYER, Sheldon, R.** et al.; Fliesler Dubb Meyer
& Lovejoy LLP, Four Embarcadero Center, Fourth Floor,
San Francisco, CA 94111-4156 (US).

(22) International Filing Date: 15 July 2002 (15.07.2002)

(25) Filing Language: English

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU,
AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU,
CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH,
GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC,
LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW,
MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SD, SE, SG,
SI, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VN,
YU, ZA, ZM, ZW.

(26) Publication Language: English

(30) Priority Data:
60/305,986 16 July 2001 (16.07.2001) US
60/7305,978 16 July 2001 (16.07.2001) US
09/975,590 11 October 2001 (11.10.2001) US
09/975,587 11 October 2001 (11.10.2001) US

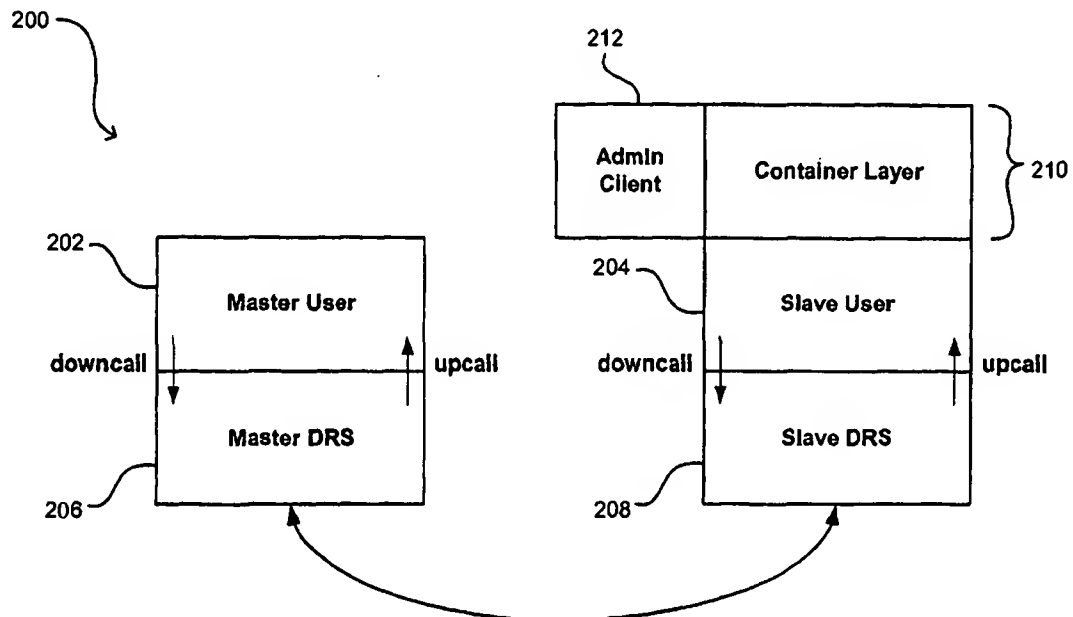
(84) Designated States (*regional*): ARIPO patent (GH, GM,
KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW),
Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM),
European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE,
ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, SK,
TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ,
GW, ML, MR, NE, SN, TD, TG).

(71) Applicant: **BEA SYSTEMS, INC.** [US/US]; 2315 North
First Street, San Jose, CA 95131 (US).

(72) Inventors: **JACOBS, Dean, Bernard**; 1747 Madera
Street, Berkeley, CA 94707 (US). **KRAMER, Reto**;
411 Green Street, #2A, San Francisco, CA 94133 (US).

[Continued on next page]

(54) Title: DATA REPLICATION PROTOCOL



(57) Abstract: Data can be replicated over a network using a one or two phase method. For the one phase method, a master server (206) containing an original copy of the data sends a version number for the current state of the data to each slave (208) on the network so that each slave can request a delta from the master. The delta that is requested contains the data necessary to update the slave to the appropriate version of the data. For the two phase method, the master server sends a packet of information to each slave. The packet of information can be committed by the slaves if each slave is able to process the commit.

WO 03/009092 A3



Published:

- with international search report
- before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(88) Date of publication of the international search report:
10 April 2003

INTERNATIONAL SEARCH REPORT

International application No.

PCT/US02/22366

A. CLASSIFICATION OF SUBJECT MATTER

IPC(7) : G06F 15/16; G06F 7/00, 17/30

US CL : 709/203; 707/2, 8, 10, 101

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 709/203, 208, 209, 242, 246

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
Please See Continuation Sheet

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 6,088,694 A (BURNS et al.) 11 JULY 2000 (11.07.2000), col. 4 line 30 to col. 6 line 43.	1-93
A	US 5,920,867 A (VAN HUBEN et al.) 06 JULY 1999 (06.07.1999), col. 6 line 50 to col. 8 line 3.	1-93
A, P	US 6,263,372 B1 (HOGAN et al.) 17 JULY 2001, (17.07.2001), col. 3 line 23 to col. 9 line 45.	1-93

☐ Further documents are listed in the continuation of Box C.

☐ See patent family annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search

18 November 2002 (18.11.2002)

Date of mailing of the international search report

28 FEB 2003

Name and mailing address of the ISA/US

Commissioner of Patents and Trademarks
Box PCT
Washington, D.C. 20231

Facsimile No. (703)305-3230

Authorized officer

Ayaz R Sheikh

Telephone No. 703 305 3900

INTERNATIONAL SEARCH REPORT

PCT/US02/22366

Continuation of B. FIELDS SEARCHED Item 3:

West, Derwent, EPO, JPO

master, slave, client, server, replicating, duplicating, copying, updating, changing, modifying, delta, phase, version